

# Analyzing Amazon's Co-Purchasing Network With Apache Spark

## Team Members

Rishabh Gautam 20bcs112

Deepesh mishra 20bcs043

Mahendra Puniya 20bcs082

RaviKath 20bcs110

## *Under Guidance of*

Dr. Animesh Chaturvedi, Assistant Professor

## Abstraction:

This analysis delves into the Amazon Product Co-purchasing Network, a dataset collected in June 2003, to unravel patterns of customer behavior and product relationships on the Amazon platform. Leveraging Apache Spark and the Connected Components algorithm, this study identifies clusters of interconnected products, shedding light on co-purchasing behavior. The findings have profound implications for enhancing recommendation systems and optimizing marketing strategies. The report underscores the significance of Apache Spark in efficiently processing and analyzing large-scale network data, offering insights that empower data-driven decision-making for e-commerce applications.

computing system. The dataset, collected in June 2003, was harvested by crawling the Amazon website and is rooted in the "Customers Who Bought This Item Also Bought" feature, a cornerstone of e-commerce recommendation systems.

The Amazon Product Co-purchasing Network provides a wealth of insights into customer interactions with products on the Amazon platform. By leveraging Apache Spark, we unlock the capability to process and analyze this vast and interconnected dataset with ease. Our analysis, powered by the Connected Components algorithm in Apache Spark's GraphX library, allows us to uncover patterns and relationships within the network, enabling us to understand how products are co-purchased, clustered, and potentially recommended together

## Introduction:

In the landscape of big data and distributed computing, Apache Spark has emerged as a powerful framework for processing and analyzing vast datasets efficiently. This report centers on the analysis of the Amazon Product Co-purchasing Network using Apache Spark, an open-source, lightning-fast, and cluster

## Dataset information:

The dataset under consideration is the Amazon Product Co-purchasing Network, collected in June 2003. This network was obtained through web crawling of the Amazon website, relying on the "Customers Who Bought This Item Also Bought" feature. In this network, if a product 'i' is frequently co-purchased with product 'j', a

directed edge is created from 'i' to 'j'. The dataset provides essential statistics, offering valuable insights into its scale and structure:

- **Nodes: 403,394**
- **Edges: 3,387,388**
- **Nodes in Largest Weakly Connected Component (WCC): 403,364 (100% of nodes)**
- **Edges in Largest Weakly Connected Component (WCC): 3,387,224 (100% of edges)**
- **Nodes in Largest Strongly Connected Component (SCC): 395,234 (98% of nodes)**
- **Edges in Largest Strongly Connected Component (SCC): 3,301,092 (97.5% of edges)**
- **Average Clustering Coefficient: 0.4177**
- **Number of Triangles: 3,986,507**
- **Fraction of Closed Triangles: 0.06206**
- **Diameter (Longest Shortest Path): 21**
- **90-Percentile Effective Diameter: 7.6**

These statistics provide a comprehensive overview of the dataset's scale, connectivity, and inherent structure, which is essential for conducting meaningful network analysis and drawing actionable insights from the data.

## Work Flow:

**Data Acquisition:** Obtain the Amazon Product Co-purchasing Network dataset, collected in June 2003, either from the source provided or the relevant data repository.

**Environment Setup:** Set up your development environment with Apache Spark and relevant libraries, ensuring that your infrastructure can handle the size and complexity of the dataset.

**Data Loading:** Load the dataset into your Spark environment using the GraphX library. Ensure that the data is appropriately structured for analysis.

**Connected Components Analysis:** Utilize the Connected Components algorithm from GraphX to identify connected components

within the network. This step will label each component with the ID of its lowest-numbered vertex.

**Component Size Calculation:** Calculate the size of each connected component by mapping the component ID to 1 and then reducing by summing up the sizes.

**Sorting Component:** Sort the connected components by size in descending order to identify the largest and potentially more significant clusters within the network.

**Result Presentation:** Present the results of the analysis, highlighting the connected components, their sizes, and any insights gained from this process.

**Interpretation and implication:** Discuss the implications of the identified connected components, exploring potential applications for improving recommendation systems, marketing strategies, and customer behavior analysis.

**Discussion of Apache Spark:** Emphasize the role of Apache Spark in efficiently processing and analyzing large-scale network data, demonstrating its capacity for distributed computing

**Conclusion:** Summarize the key findings and the significance of the analysis, highlighting the value of modern data science tools in unlocking actionable insights from extensive datasets.

**Report Generation:** Compile the findings and insights into a comprehensive report, including an abstract, introduction, dataset information, objectives, analysis, and a conclusion.

**Further Analysis and Application:** Consider opportunities for additional analysis or applications based on the specific objectives and needs of your project. This workflow guides you through the process of leveraging Apache Spark to analyze the Amazon Product Co-purchasing Network, offering insights into customer behavior and product relationships on the Amazon platform.

## RESULT :

The largest connected components within the network play a critical role in understanding customer behavior and product relationships. These findings have practical implications for improving recommendation systems, marketing strategies, and customer behavior analysis, particularly in the context of e-commerce.

### Conclusion:

demonstrated its prowess in processing and analyzing large-scale network data efficiently. Its distributed computing capabilities allowed for the exploration of intricate patterns and relationships within the dataset.

This analysis underscores the power of modern data science tools in extracting actionable insights from extensive datasets, and it highlights the potential of Apache Spark in unlocking the value of large-scale network data.

1. Leskovec, J., Adamic, L., & Adamic, B. (2007). The Dynamics of Viral Marketing. *ACM Transactions on the Web (ACM TWEB)*, 1(1).
2. Apache Spark - Official Website. (<https://spark.apache.org/>)
3. GraphX - Apache Spark. (<https://spark.apache.org/graphx/>)